Let's return to a topic we had in MATH 300.

$$X_1, \ldots, X_n \overset{iid}{\sim} \rho \, N(\theta_0, \sigma_0^2) + (1-\rho) N(\theta_1, \sigma_1^2) = \sum_{m=1}^{M} \rho_m N(\theta_m, \sigma_m^2)$$

$$\text{where } M=2$$

This was called a mixture model. Parameters? $\{\rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2\}$

$$f(\vec{X}; \rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2) = \mathcal{L}(\rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2 | \vec{X}) = \prod_{i=1}^{n} \rho \frac{1}{\sqrt{2\pi\sigma_0^2}} \cdots \cdots + (1-\rho) \cdots$$

MLE's?

$$\ell(\rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2 | \vec{X}) = \sum_{i=1}^{n} \ln \left( \rho \cdots \right.$$

Very difficult to take derivatives and set $=0$ to solve. As $M$ increases, even more difficult.

Enter "data augmentation" (1987)

What if we knew

$$X_1 \text{ belongs to } N(\theta_1, \sigma_1^2), \quad X_2 \text{ belongs to } N(\theta_2, \sigma_2^2), \quad X_3 \ldots ?$$

let $I_i = \mathbb{1}$ $i^{th}$ observation "comes from" $N(\theta_0, \sigma_0^2)$

$$\Rightarrow I_1, \ldots, I_n \overset{iid}{\sim} \text{Bern}(\rho) \qquad \text{"data augmentation" "add" "new data", the } I_i\text{'s.}$$

$$f(\vec{X}, \vec{I} | \rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2) = f(\vec{X} | \rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \vec{I}) \, P(\vec{I} | \rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2)$$

$$= \prod_{i=1}^{n} \left( \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(X_i - \theta_0)^2} \right)^{I_i} \left( \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(X_i - \theta_1)} \right)^{1-I_i} \prod_{i=1}^{n} \rho^{I_i} (1-\rho)^{1-I_i}$$

Hw: verify $f(\vec{x} \mid \cdots) = \sum\limits_{S_{\vec{I}}} f(\vec{x}, \vec{I} \mid \theta)$

How does this help? We can now find MLE's!

$$\ell(\rho, \theta_0, \sigma_0^2, \theta_1, \sigma_1^2; \vec{x}, \vec{I}) = \sum_{i=1}^{n} \ln\left( \cdots \qquad \right)$$

$$= \sum_{i=1}^{n} I_i \left( -\tfrac{1}{2}\ln(2\pi) - \tfrac{1}{2}\ln(\sigma_0^2) - \tfrac{1}{2\sigma_0^2}(x_i - \theta_0)^2 \right)$$

$$+ (1-I_i)\left( -\tfrac{1}{2}\ln 2\pi - \tfrac{1}{2}\ln(\sigma_1^2) - \tfrac{1}{2\sigma_1^2}(x_i - \theta_1)^2 \right)$$

$$+ I_i \ln(\rho) + (1-I_i)\ln(1-\rho)$$

$$= -\tfrac{1}{2}\ln(2\pi)\sum I_i - \tfrac{1}{2}\ln(\sigma_0^2)\sum I_i - \tfrac{1}{2\sigma_0^2}\sum(x_i-\theta_0)^2 I_i$$

$$-\tfrac{1}{2}\ln(2\pi)\sum(1-I_i) - \tfrac{1}{2}\ln(\sigma_1^2)\sum(1-I_i) - \tfrac{1}{2\sigma_1^2}\sum(x_i-\theta_1)(1-I_i)$$

$$+ \ln(\rho)\sum I_i + \ln(1-\rho)\sum(1-I_i)$$

let $n_0 = \sum I_i$
let $n_1 = \sum 1-I_i$

$$\frac{\partial \ell}{\partial \theta_0} = -\frac{1}{2\sigma_0^2}\sum \frac{\partial}{\partial\theta_0}(x_i-\theta_0)^2 I_i = -\frac{1}{2\sigma_0^2}\frac{\partial}{\partial\theta_0}\left[\sum x_i^2 I_i - 2\theta_0\sum x_i I_i + \theta_0^2 \sum I_i\right] \overset{set}{=} 0$$

$$\Rightarrow -\sum x_i I_i + \theta_0 \sum I_i = 0 \Rightarrow \hat\theta_0^{MLE} = \frac{\sum x_i I_i}{n_0}, \quad \hat\theta_1^{MLE} = \frac{\sum x_i(1-I_i)}{n_1} \quad \underline{\underline{Hw}}$$

$$\frac{\partial \ell}{\partial \sigma_0^2} = -\frac{\sum I_i}{2\sigma_0^2} + \frac{\sum(x_i-\theta_0)^2 I_i}{2(\sigma_0^2)^2} \overset{set}{=} 0 \Rightarrow -\sum I_i + \frac{\sum(x_i-\theta_0)^2 I_i}{\sigma_0^2} = 0 \Rightarrow \hat\sigma_0^{2\,MLE} = \frac{\sum(x_i-\theta_0)^2 I_i}{n_0}$$

$$\hat\sigma_1^{2\,MLE} = \frac{\sum(x_i-\theta_0)^2(1-I_i)}{n_1}$$

$$\frac{\partial \ell}{\partial \rho} = \frac{n_0}{\rho} - \frac{n_1}{1-\rho} \overset{set}{=} 0 \Rightarrow n_0(1-\rho) = n_1\rho \Rightarrow n_0 - n_0\rho = n_1\rho \Rightarrow n_0 = (n_0+n_1)\rho$$

$$\Rightarrow \hat\rho^{MLE} = \frac{n_0}{n_0+n_1}$$

Okay great... who cares? I don't know $I_1 ... I_n$!
So all of this is useless, but a fun math exercise!

What if we can estimate the $I_i$'s?

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ since $I_i$ is Bernoulli

let $\hat{\hat{I}}_i := E\left[I_i \mid \theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \rho, X_i\right] = P(\overset{A}{I_i=1} \mid \overbrace{\theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \rho}^{C}, \overset{B}{X_i})$

$$P(A|C) = \frac{P(B|A)P(A)}{P(C)}$$

$$P(A|BC) = \frac{P(B|A,C)\, P(A|C)}{P(B|C)}$$

$$= \frac{f(X_i, I_i=1 \mid \theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \rho)}{f(\qquad\qquad)}$$

$$= \frac{f(X_i \mid \theta_0, \sigma_0^2, \theta_1, \sigma_1^2, I_i=1)\; \overbrace{P(I_i=1 \mid \theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \rho)}^{\rho}}{f(X_i \mid \theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \rho)}$$

$$= \frac{\rho \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(X_i - \theta_0)^2}}{\rho \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(X_i-\theta_0)^2} + (1-\rho)\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(X_i-\theta_1)^2}}$$

Why not do the following:

Step 0: Initialize $\vec{\theta}_0 = 0, \sigma_0^2 = 1, \vec{\theta}_1 = 1, \sigma_1^2 = 1, \rho = 0.5$

Step 2: Compute $\vec{\theta}_0, \vec{\sigma}_0^2, \vec{\theta}_1, \vec{\sigma}_1^2, \hat{\rho}$ MLE based on $\hat{I}_1, ..., \hat{I}_n$

Step 1: Compute $\hat{I}_1, ..., \hat{I}_n$ based on $\vec{\theta}_0, \vec{\sigma}_0^2, \vec{\theta}_1, \vec{\sigma}_1^2, \hat{\rho}$

Step 3: Repeat steps 1, 2 until $\|\vec{\theta}_t - \vec{\theta}_{t-1}\| < \varepsilon$

This is called the Expectation-Maximization Algorithm $\left(EM, 1977\right)$

$\qquad\qquad\qquad\uparrow\qquad\qquad\qquad\uparrow$
$\qquad\qquad\quad$ step 1 $\qquad\qquad$ step 2

This is prone to converge to $\hat{\vec{\theta}}^{MLE}$. Also, it's possible to derive the Fisher Information estimates. This gives you asymptotically valid inference via the MLE normal thm.

The EM algorithm is very general. Not just for mixture models inference.

Note: easy to upgrade to $M>2$ $\quad \vec{I}_1,...,\vec{I}_n \overset{iid}{\sim} Multinom(1, \vec{e})$ $\quad$ [MA HW]

Bayesian Approach. Use data augmentation as well. Treat $\vec{I}$ as parameter

Let $\quad f(\theta_0) \propto 1, \quad f(\sigma_0^2) \propto \frac{1}{\sigma_0^2}, \quad f(\theta_1) \propto 1, \quad f(\sigma_1^2) \propto \frac{1}{\sigma_1^2}, \quad f(\rho) = U(\cdot,1)$

$\qquad\qquad\qquad \underbrace{\phantom{xxxxx}}_{Jeffrey's} \qquad\qquad\qquad \underbrace{\phantom{xxxxx}}_{Jeffrey's} \qquad\qquad \underset{Laplace}{\phantom{x}}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad f(\vec{I}) \propto 1$

$$f\left(\theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \rho, \vec{I} \mid \vec{X}\right) \propto (\sigma_0^2)^{-n/2-1} e^{-\frac{1}{2\sigma_0^2}\sum(X_i-\theta_0)^2 I_i}$$

$$(\sigma_1^2)^{-n/2-1} e^{-\frac{1}{2\sigma_1^2}\sum(X_i-\theta_1)^2(1-I_i)}$$

$$e^{n_0} (1-\rho)^{n_1}$$

Try to make a Gibbs Sampler

$$f(\theta_0 \mid \text{---}) \propto e^{-\frac{1}{2\sigma_0^2}\left(\sum X_i^2 I_i - 2\theta_0 \sum X_i I_i + \theta_0^2 \sum I_i\right)}$$

$$\propto e^{\frac{\sum X_i I_i}{\sigma_0^2}\theta_0 - \frac{n_0}{2\sigma_0^2}\theta_0^2} = e^{a\theta_0 - b\theta_0^2}$$

$$\propto N\left(\frac{\sum X_i I_i}{n_0}, \frac{\sigma_0^2}{n_0}\right)$$

$$f(\theta_1 \mid \text{---}) \propto N\left(\frac{\sum X_i(1-I_i)}{n_1}, \frac{\sigma_1^2}{n_1}\right)$$

$$f(\sigma_0^2 \mid -\!\!\!-) \propto (\sigma^2)^{-\frac{4_0}{2}-1} e^{-\frac{\xi(x_i - \theta_0)^2 I_i / 2}{\sigma_0^2}}$$

$$\propto \text{InvGamma}\left(\frac{4_0}{2}, \frac{\xi(x_i - \theta_0)^2 I_i}{2}\right)$$

$$f(\sigma_1^2 \mid -\!\!\!-) \propto \text{InvGamma}\left(\frac{4_1}{2}, \frac{\xi(x_i - \theta_0)^2 (1 - I_i)}{2}\right)$$

$$f(\rho \mid -\!\!\!-) \propto \text{Beta}(n_0 + 1, n_1 + 1)$$

$$P(I_i \mid -\!\!\!-) \propto e^{-\frac{1}{2\sigma^2}(x_i - \theta_0) I_i} \, e^{-\frac{1}{2\sigma_1}(x_i - \theta_1)(1 - I_i)} \, \rho^{I_i} (1 - \rho)^{1 - I_i}$$

$$= \left(\rho \, e^{-\frac{x_i - \theta_0}{2\sigma^2}}\right)^{I_i} \left((1 - \rho) \, e^{-\frac{x_i - \theta_1}{2\sigma_1}}\right)^{1 - I_i}$$

$$\propto \text{Bern}\left(\frac{\rho \, e^{-\frac{x_i - \theta_0}{2\sigma_0^2}}}{\rho \, e^{-\frac{x_i - \theta_0}{2\sigma_0}} + (1 - \rho) \, e^{-\frac{x_i - \theta_1}{2\sigma_1}}}\right)$$

We made a Gibbs Sampler!

---

Remember, Stan does not allow for discrete params.

So in Stan, we write the lik without down augmentation